# Computational identification of discriminating features of pathogenic and symbiotic type III secreted effector proteins

Koji Yahara[†]    Ying Jiang[†]    Takashi Yanagawa[††]

Type III secr etion sy stems (T3SS) de liver bac terial p roteins, o r "e ffectors", into eukaryotic host cells , inducing ph ysiological respo nses in the hosts. E ffector proteins have been considered virulence factors of pathogenic bacteria, but T3SSs have now been found in sy mbiotic bacteria as well. Whether an y phy sicochemical dif ference exists between th e two ty pes of ef fectors remains unkn own. In th is work, we combined computational statis tical and machine-learning methods to find th e ph ysicochemical differences. The most dis criminating set of fe atures in a d ataset of ph ysicochemical features was de termined us ing g eneralized B ayesian infor mation c riteria and k ernel logistic reg ression. Classification per formance was ex amined u sing a suppo rt vecto r machine. Int erdependence among the most discriminating featu res was explored b y graphical modeling, and th e most dis criminating region was inv estigated by s liding window analysis.

# 病原細菌と共生細菌の III 型分泌装置のエフェクタータンパク質を区別する特徴は何か？

矢原耕史[†]    姜英[†]    柳川堯[††]

近年、細菌がホスト細胞に送り込むエフェクタータンパク質が、病原細菌と共生細菌の双方に存在することが注目されているが、その差異は未だ明らかでない。本研究では、両者の差異を司る物理化学的特徴セットをカーネルロジスティック回帰の情報量基準によって抽出し、その判別性能を SVM によって評価し、さらにその相互依存関係及び最も特徴的な領域を、グラフィカルモデリングと sliding window 解析によって明らかにした。

## 1. Introduction

Type III secr etion sy stems ( T3SS) are co mplex secretion machine s that deliver bacterial protei ns called ef fectors i nto euka ryotic ho st cells t hrough an injec tisome du ring infection ( 1, 2) . T3SS-secreted ef fector prot eins in duce ph ysiological resp onses in t heir hosts, suc h as cy toskeletal rearrange ment to promote bacteri al attachment and invasi on, interference wit h cellular traf ficking proce sses, cy totoxicity (2), induction o f apopto sis o f macrophages ( 3), di sruption of tight juncti ons (4), and m icrotubule de stabilization (5 ). These ef fector protein functi ons are consi dered causes of vi rulence in pathogenic bacteria such a s *Yersinia* s pecies (s pp.), *Chlamydia* s pp., *Salmonella* s pp., *Shigella* s pp., and enteropathogenic *Escherchia coli*. However, T3SSs are also found in sy mbiotic bacteria (6, 7), and a genome analy sis of a Chlamy dia-related symbiont o f free-living amoebae sugge sts that the origins of T3SSs may be unrelated to virulence (8).

Common features of T3SS effector proteins in pathogenic and symbiotic bacteria can be identified by computational methods (9, 10 ). While T3SS effector protein s were origi nally not tho ught to share any comm on feature s (11), recent studies u sing machine-learning approaches have identified comm onalities in the N-terminus of ef fectors, mainly in am ino acid composition. One st udy (9) analy zed bot h pathogenic and sy mbiotic T3SS ef fector proteins, and f ound a si gnature in the N-terminus that i s taxono mically univer sal an d conserved.

The sy mbiotic T3SS ef fector proteins, ho wever, have different functio ns t han th e pathogenic ef fectors. Sy mbiotic ef fectors o f rhizobia, for example, modulate host-plant reactions, that lead to the formation of functional nodules (12, 13). Putative effector proteins of the tset se fl y endosy mbiont, *Sodalis glo ssinidius*, specifically facilitate the host cell cytoskeletal rearrange ments necessary for ba cterial entr y, although the number of gene s encoding effector proteins is smaller in the s ymbiotic regions than in the homol ogous islands in pathogenic bacteria (14). Ho mologs of the sy mbiotic region s are also fo und i n endosymbionts of grain weevil s, *Sitophilus o ryzae* and *S. zea mais,* in which T3SS genes are suggested to fu nction durin g a specific stage of weevil developme nt (14). Even if the signature amin o acid sequence in the N-te rminus is con served amo ng p athogenic and symbiotic T3SS ef fector p roteins, the se fu nctional dif ferences exist. We were interested in

[†]Division of Biostatistics, Kurume University School of Medicine
　久留米大学医学研究科バイオ統計学群
[††]Biostatistics Center, Kurume University
　久留米大学バイオ統計センター

finding t he phy sicochemical di fferences betwe en pathogenic   and sy mbiotic  T3SS ef fector proteins that might be responsible for these functional differences.

In this work, we com bined com putational st atistical and mac hine-learning approache s to addre ss thi s issue.   From a  dataset of physicochemical features prepared from pathogenic and symbiotic T3SS effector proteins, the most discriminating set of feature s was determined using generalized Baysian information criteria and kernel logi stic re gression.   Classification performance using the i dentified di scriminating feature s was examined u sing su pport vect or machine (SVM ).   The res ults clearly s howed differences in am ino acid compos ition.   The most discriminating set of seven features were identified and successfully used to classify the effectors, with a sensitivity and specificity of over 80%.   In addition, interdependence among the  most discri minating seven   features wa s re vealed by  graphical modeli ng. The   most discriminating r egion for the  most di scriminating seven featur es was deter mined by  sliding window analysis.

## 2.   Materials and Methods

### 2.1   Dataset

We collected the 57 currentl    y available am ino acid s equences of s ymbiotic T3S S effector proteins from the literature (9, 15), and the same number of amino acid sequences for pathogenic T3SS effector proteins (9).

For each e ffector protei n am ino acid s equence, we calculated the phy   sicochemical features, 41 in  total, of cha rge, isoetectric p oint, num ber of  proteolytic enzy me or reagent cleavage s ites, mole percenta ge of each a  mino acid and a   mino acid groups defined i   n EMBOSS (16), and signal pepti de probability.   The list of 41 physicochemical features used in this st udy is i n Table 1. Signal peptide pr obability was calculated b y SignalP 3.0 (17), a nd others feature s  were calculated by  EM BOSS (16).   These  were used  as att ributes in o ur classification analysis.

### 2.2   Feature selection

We first used t he Lepage test   for the lo cation-dispersion di fference between the two groups (18).   The top 10 discriminating features were chosen by the order of their p-values in the test statistics.   The p-values of all of these candidate features were less than 0.001.

For t hese candidate features  , we exa mined a ll combinations, $2^{10}-1$, a s explanatory variables in the kernel logistic regression (KLR), which is one of the kernel-learning methods suitable for binary-pattern recognition problems (19, 20).   Let $y_i$ be a binary observed

Table 1. Biochemical features used as attributes of effector proteins

| No. | Description |
|---|---|
| 1 | Number of potentially antigenic regions of a protein sequence[1] |
| 2 | Number of proteolytic enzyme or reagent cleavage sites[1] |
| 3 | Number of secondary structure[1] |
| 4 | Hydrophobic moment[1] |
| 5 | Average residue weight[1] |
| 6 | Charge[1] |
| 7 | Isoelectric point[1] |
| 8 | Molar extinction coefficient[1] |
| 9 | Extinction coefficient at 1 mg/ml[1] |
| 10 | Probability of protein expression in E. coli inclusion bodies[1] |
| 11-30 | Mole percentage of each amino acid[1] 11:Ala, 12:Cys, 13:Asp, 14:Glu, 15:Phe, 16:Gly, 17:His, 18:Ile, 19:Lys, 20:Leu, 21:Met, 22:Asn, 23:Pro, 24:Gln, 25:Arg, 26:Ser, 27:Thr, 28:Val, 29:Trp, 30:Tyr |
| 31 | Mole percentage of tiny amino acids[1] (A+C+G+S+T) |
| 32 | Mole percentage of small amino acids[1] (A+B+C+D+G+N+P+S+T+V) |
| 33 | Mole percentage of aliphatic amino acids[1] (A+I+L+V) |
| 34 | Mole percentage of aromatic amino acids[1] (F+H+W+Y) |
| 35 | Mole percentage of non-polar amino acids[1] (A+C+F+G+I+L+M+P+V+W+Y) |
| 36 | Mole percentage of polar amino acids[1] (D+E+H+K+N+Q+R+S+T+Z) |
| 37 | Mole percentage of charged amino acids[1] (B+D+E+H+K+R+Z) |
| 38 | Mole percentage of basic amino acids[1] (H+K+R) |
| 39 | Mole percentage of acidic amino acids[1] (B+D+E+Z) |
| 40 | Number of clea vage sites be tween s ignal s equence and m ature exported protein[1] |
| 41 | Signal peptide probability[2] |

[1] calculated by EMBOSS (16).   [2] calculated by SignalP (17).

variable and $p(\mathbf{x}_i)$ be its co nditional distribution given $\mathbf{x}_i$, the n the likelihood function was given by

$$L = \prod_{i=1}^{n} p(\mathbf{x}_i)^{y_i} (1 - p(\mathbf{x}_i))^{1-y_i} \tag{1}$$

and log-likelihood function became

$$\begin{aligned}\log L \\ = \sum_{i=1}^{n} y_i \log \frac{p(\mathbf{x}_i)}{1 - p(\mathbf{x}_i)} + \log(1 - p(\mathbf{x}_i))\end{aligned} \tag{2}$$

in whic h the unknown quantity $p(\mathbf{x}_i)$ was modeled using the ra dial basi s ker nel fu nction $K(\mathbf{x}_j, \mathbf{x}_i)$ as

$$f(\mathbf{x}_i) = \log \frac{p(\mathbf{x}_i)}{1 - p(\mathbf{x}_i)} = \sum_{j=0}^{n} \alpha_j K(\mathbf{x}_j, \mathbf{x}_i) \tag{3}$$

where

$$\mathbf{K}_{ij} = K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\sigma \| \mathbf{x}_i - \mathbf{x}_j \|^2) \tag{4}$$

and $\sigma$ is the kernel param eter. The soluti on of the para meter vector $\hat{\boldsymbol{\alpha}}$ was calculated using the following penalized log-likelihood function

$$\frac{1}{n}\{\sum_{i=1}^{n} y_i f(\mathbf{x}_i) - \log[1 + \exp(f(\mathbf{x}_i))]\} - \frac{\lambda}{2} \boldsymbol{\alpha}^T \mathbf{R} \boldsymbol{\alpha} \tag{5}$$

where

$$R = \begin{pmatrix} 1 & \mathbf{1}_n^T \\ \mathbf{1}_n & \mathbf{K} \end{pmatrix}$$

by Fisher's scoring methods.

To s elect the b est combination of the 10 candidate features , we us ed a generalized Bayesian infor mation crite rion (GB IC) ( 21). Using the likelihood f unction $L(\boldsymbol{\alpha})$ in equation (1) and the multivari ate norm al prior density $\pi(\boldsymbol{\alpha} | \lambda)$ for th e parameter ve ctor $\boldsymbol{\alpha}$ defined by

$$\pi(\boldsymbol{\alpha} | \lambda) = (2\pi)^{-r/2} (n\lambda)^{r/2} | R |_+^{1/2} \exp(-\frac{n\lambda}{2} \boldsymbol{\alpha}^T R \boldsymbol{\alpha}) \tag{6}$$

GBIC was defined as

$$GBIC = -2 \log \int L(\boldsymbol{\alpha}) \pi(\boldsymbol{\alpha} | \lambda) d\boldsymbol{\alpha} \tag{7}$$

and $R$ was the same as that of equation (5), $r$ was the rank of $R$, and $| R |_+$ was the product of $r$ nonzer o eigenvalues of $R$. O nce $\hat{\boldsymbol{\alpha}}$ was o btained, GBI C was calculated through th e Laplace approximation

$$\begin{aligned}&-2 \log \int \exp(nl_\lambda(\boldsymbol{\alpha})) d\boldsymbol{\alpha} \\ &= -2 \log \{ \frac{(2\pi / n)^{(n+1)/2}}{| J_\lambda(\hat{\mathbf{a}}) |^{1/2}} \exp(nl_\lambda(\hat{\mathbf{a}})) \} \{1 + O(n^{-1})\}\end{aligned} \tag{8}$$

where

$$l_\lambda(\boldsymbol{\alpha}) = \frac{1}{n} \log L(\boldsymbol{\alpha}) + \frac{1}{n} \log \pi(\boldsymbol{\alpha} | \lambda)$$

$$J_\lambda(\hat{\boldsymbol{\alpha}}) = -\frac{\partial^2 l_\lambda(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha} \partial \boldsymbol{\alpha}^T},$$

GBIC was computed for each com bination of 10 features , and the com bination with the minimum GBI C wa s deter mined as ex planatory variable of KLR. During the feature selection, values of kernel parameter $\sigma$ and hyper parameter $\lambda$ were given in the range of 1E-3 to 1E+3 ( $\sigma$ ) or to 1E+4 ( $\lambda$ ) for each set of explanatory features.

### 2.3 Classification performance

Classification performance u sing di scriminating feature s ide ntified by GBI C of KRL was analy zed using SVM based on the approx imate relationship between KRL and the SVM (19). To determ ine t he adva ntage of t he most di scriminating feature s, a misclassification rate was eval uated by le ave-one-out cr oss-validation fo r each combi nation of $k$-features t hat attained the minimum GBIC in $_{10}C_k$ combinations ($k$=1,…,10). The r esults are summarized Figure 2, which illustrates the misclassification rates, with the num ber of features on the horizontal axis. We u sed sv m functio n of e1071 package (E. Dimitriad ou, K. Hornik, F . Leisch, D. Meyer, and A. Weingessel) in R.

### 2.4 Graphical modeling

To explore interdependence among the mo st discriminating features identified by GBIC of KLR, we u sed graphical modeling develo ped b y Imoto et al. (22, 23) which combine s non-linear n onparametric regr ession with ra dial basis and Bay esian networ k, and was originally developed for esti mating genetic ne tworks and fu nctional relatio nships bet ween genes. Non-l inear nonparametric regression enabled u s to capture directed dependencie s among the featu res without ad vance knowledge about their rela tionships. Bayesian network is a powe rful, graph-theoretic approach f or expressing i nterdependence amon g variables a s networks.

Calculations were conducted by MA TLAB R2008b (The Mathwork s Inc. ) based on NETLAB (24), the Ba yes net toolbox (BNT) for Matlab (25), and BNT structure learning package (26).

### 2.5 Sliding window analysis

N-terminal re gions from the 1 st to 97 th re sidue were analy zed, with the window size varying from 8-50, and the s tarting pos ition varying from 1 to 50. For each window , a dataset of the most di scriminating feature s wa s created, and clas sification wa s conducted by SVM.

## 3. Results

### 3.1 Identification of discriminating features

A plot of mini mum GBIC for $_{10}C_k$ co mbination of feature s u sed in KL R wa s gi ven in Figure 1 taking the numbe r of features, $k,$ on the horizontal axis . The figure s hows that the minimum GBIC tends to decrease as the numbe r of features inc rease, take the sm allest value when the number of features is seven, and incr ease at greater than seven features. The seven features that att ained the sm allest m inimum GBIC we re as fol lows: average res idual wei ght, mole percentage of Ala, Asp, Ile, tiny amino acids, small amino acids, and acidic amino acids.
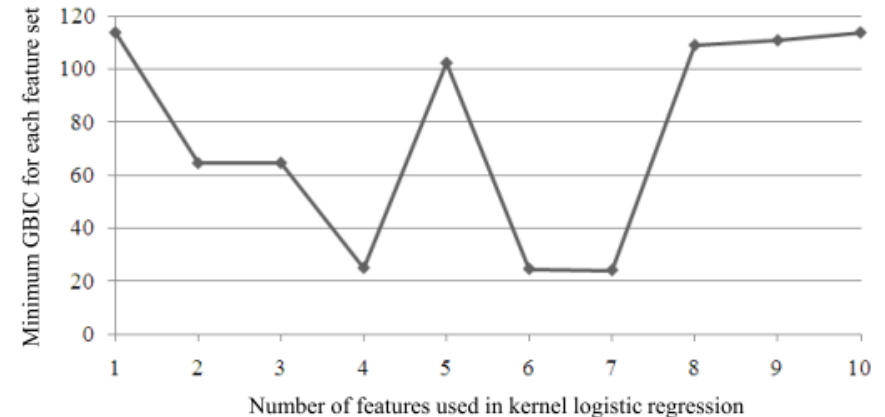


**Figure 1.** Plot of minimum GBIC against number of features used in kernel logistic regression.

### 3.2 Classification performance using the most discriminating features

Misclassification rates using the discriminating features identified by GBIC of KLR are plotted in Fi gure 2, taking the nu mber of featu res on h orizontal axis. The plot of mi nimum GBICs (Figure 1) and misclassification rat es showed parallel tendencies. The be st classification perfor mance (84. 2%) wa s obtain ed using a combination of the seven feature s that gave the sm allest m inimum GBIC (Figure 2A). T he best per formance with se ven features was ne arly identical t o the res ults obtained when all 41 features were us ed. The seven di scriminating features had a specificity of 85.5% an d a sensiti vity of 83.1% ( Figure 2B).
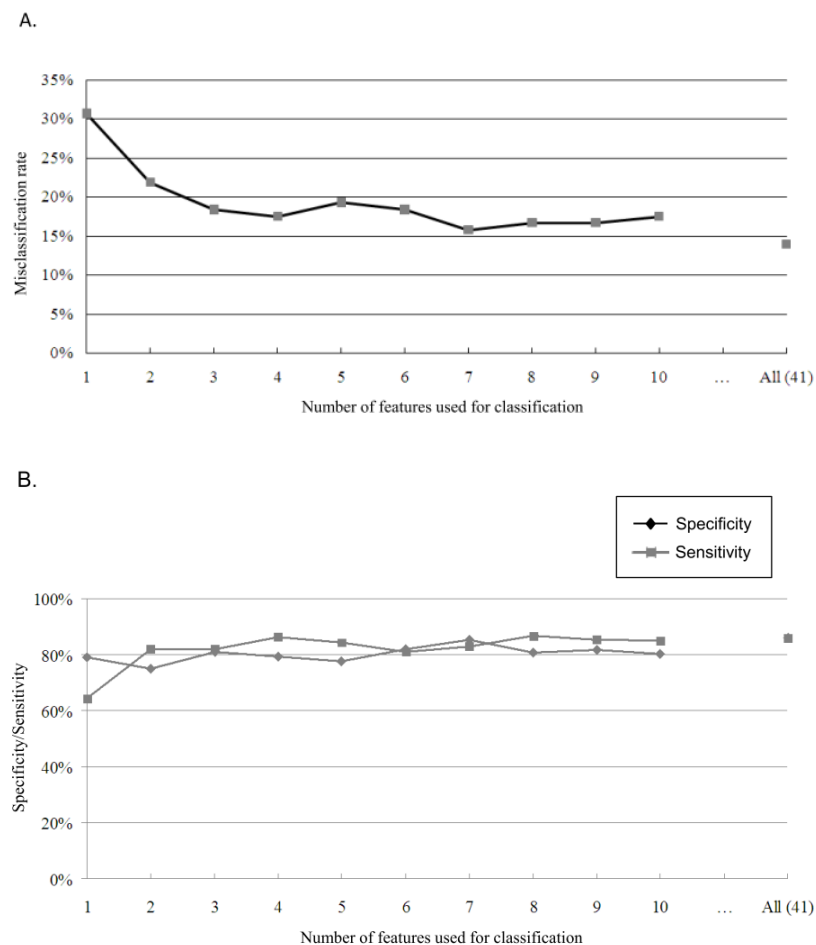
A.



B.



**Figure 2. Classification performance using the discriminating features identified by GBIC of KLR.** Misclassification rate for each com bination of k-features that attained the minimum G BIC in $_{10}C_k$ combinations (k=1,…,10). Classification using all 41 feature s was also conducted, and the misclassification ra te is at "All (41)" of the x-axis. **(A) Misclassification rate. (B) Specificity and sensitivity.**

### 3.3 Interdependence among the most discriminating features as a graph structure

The interdependence among the seven most discriminating features was represented in a directed-graph structure ( Figure 3), in whic h the mole percentage of isol eucine, and a combination of alanine and average res idue weight were po sitioned at the bottom end. The three feature s a re repre sentative of the directed-graph structure and ha ve be en selected by KLR at one or two features. Figure 2 shows that classification accuracy was about 70% for the mole percentag e of isoleucine, and nearly 80% for a com bination of alanine and average residue weight
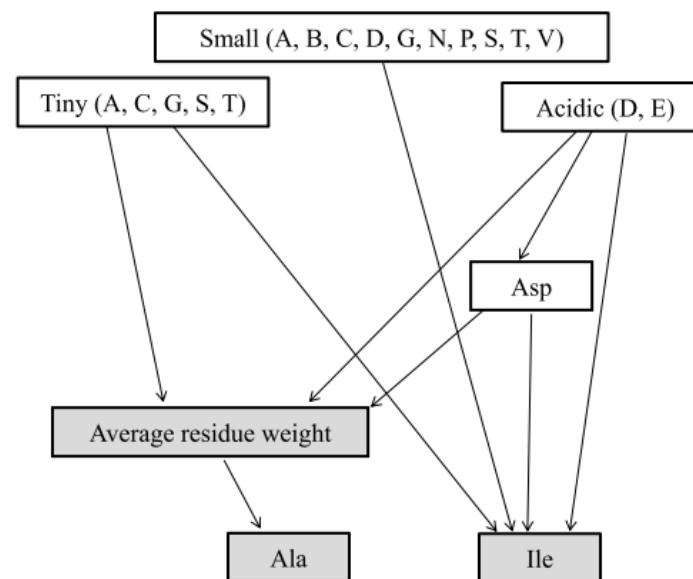


**Figure 3. Graph structure showing interdependence among the most-discriminating features.** Directed dependencies detected b y nonparametric regressi on are depicted b y arrows whos e heads indicate res ponse variable s and tails indicate explanatory variables. Colours are t he discri minating features ide ntified by GBIC, when the nu mber of feature s is one or two.

### 3.4 Identification of the most discriminating region

Sliding window analysis with variable window sizes and starting points is in Table 2. The region that gave the highest discrimination among the seven most-discriminating features was 48-95 residues from the N-terminus (N48-95), which gave a classification accuracy of 83.3% (Figure 4). Almost all of the second and third most-discriminating regions overlapped this region, supporting the hypothesis that the discriminating signature between pathogenic and symbiotic T3SS effector proteins was in this region.

**Table 2. Results of sliding window analysis**

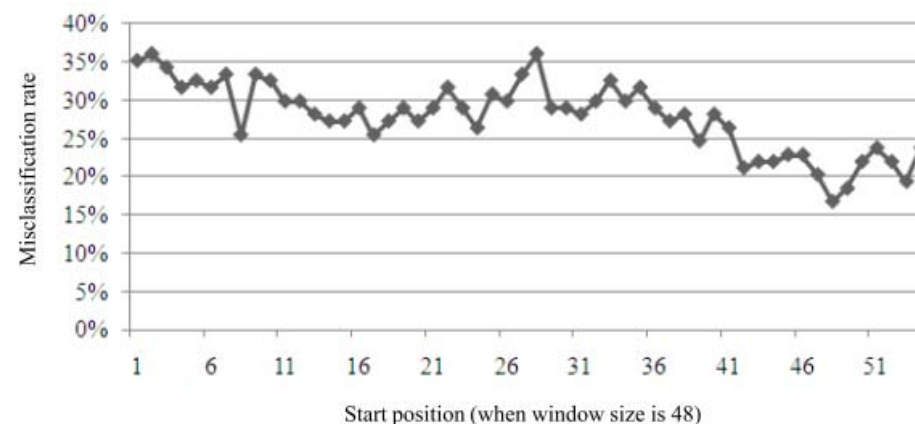| Region | Misclassification rate | Starting point | Window size |
|--------|------------------------|----------------|-------------|
| N48-95 | 0.167 | 48 | 48 |
| N49-95 | 0.175 | 49 | 47 |
| N48-93 | 0.184 | 48 | 46 |
| N48-96 | 0.184 | 48 | 49 |
| N49-89 | 0.184 | 49 | 41 |
| N49-90 | 0.184 | 49 | 42 |
| N49-96 | 0.184 | 49 | 48 |
| N9-36  | 0.184 | 9  | 28 |
| N40-89 | 0.193 | 40 | 50 |
| N47-93 | 0.193 | 47 | 47 |
| N47-96 | 0.193 | 47 | 50 |
| N48-92 | 0.193 | 48 | 45 |
| N48-94 | 0.193 | 48 | 47 |
| N49-93 | 0.193 | 49 | 45 |
| N50-96 | 0.193 | 50 | 47 |
| N65-97 | 0.193 | 65 | 33 |



**Figure 4. Plot of misclassification rate by sliding window analysis with window size 48.** As shown in Table 2, misclassification rate is lowest when the analysis start position is 48 (*i.e.* for region N48-95), and when the window size is 48, which gives the best classification performance.

### 3.5 Directions of differences of the discriminating features

The differences of the seven most discriminating features between pathogenic and symbiotic T3SS effector proteins are in Table 3, with "+" meaning "more common in symbiotic proteins". Results are given for all regions, and for the most discriminating region, N48-95. The patterns of differences were almost equivalent between all regions and the most-discriminating region, supporting the hypothesis that N48-95 was the representative region that distinguished between pathogenic and symbiotic T3SS effector proteins. By mole percentage of amino acids, isoleucine decreased in symbiotic proteins, while the other amino acids (alanine, aspartic acid, acidic amino acids, tiny amino acid, small amino acid) increased in symbiotic proteins. The tendency was found both in all regions, and in the most discriminating N48-95 region.

**Table 3. Directions of the differences of the most discriminating features**

| Pathogen | (all regions) | | symbiont (all regions) | | direction* | |
|---|---|---|---|---|---|---|
| Feature Mean | | SD | mean | SD | Mean | SD |
| Ile (Molar %) | 5.74 | 2.25 | 3.98 | 1.52 | - | - |
| Average residue weight | 109.30 | 4.12 | 108.98 | 2.80 | - | - |
| Ala (Molar %) | 8.28 | 3.07 | 10.99 | 2.80 | + | - |
| Asp (Molar %) | 4.49 | 1.91 | 16.01 | 1.60 | + | - |
| Acidic (Molar %) | 10.79 | 3.91 | 11.70 | 2.21 | + | - |
| Tiny (Molar %) | 31.58 | 6.98 | 32.94 | 3.90 | + | - |
| Small (Molar %) | 51.97 | 6.95 | 54.53 | 4.41 | + | - |
| pathogen | (N48-95) | | symbiont (N48-95) | | direction* | |
| Feature Mean | | SD | mean | SD | Mean | Region |
| Ile (Molar %) | 5.30 | 4.09 | 3.84 | 2.81 | - | - |
| Average residue weight | 109.26 | 6.41 | 109.79 | 4.34 | + | - |
| Ala (Molar %) | 9.06 | 4.55 | 10.78 | 5.25 | + | + |
| Asp (Molar %) | 3.07 | 2.27 | 5.88 | 3.46 | + | + |
| Acidic (Molar %) | 8.92 | 5.69 | 11.15 | 4.91 | + | - |
| Tiny (Molar %) | 32.35 | 10.64 | 33.08 | 8.05 | + | - |
| Small (Molar %) | 51.68 | 10.95 | 53.91 | 6.69 | + | - |

\* from pathogenic to symbiotic ("+" means "more in symbiotic proteins")

## 4. Discussion

In this wor k, we identified the seven most-discriminating featu res between pat hogenic and symbiotic T3SS effector proteins, using a large combination of phy sicochemical features, analyzed b y GBIC of KLR. The identified features were successfully us ed to class ify the proteins by SVM, with sensitivities and specificities of over 80%.

The seven mo st-discriminating features were those related to a mino acid compositi on. No ot her hi gher-order inf ormation wa s f ound to be as di scriminating by GBIC o f KL R. Interestingly, recently reported common features of T3SS ef fectors we re als o found to be amino acid com position or shared s equence m otif. E mbedded features i n the am ino aci d sequence or composition may be a characteristic of T3SS effector proteins.

The most discriminating region between pathogenic and symbiotic effector proteins was 48-95 resi dues fro m the N-ter minus. The classic signal peptide secretion signal is 15-40 residues f rom t he N-ter minus (27). Common feature s of T3SS ef fectors protein s were recently found t o be em bedded in 30 (1 0) or u p to 50 resi dues (9) at the N-ter minus. These findings are co mplementary with our s becau se the dif ferences between p athogenic an d symbiotic ef fector protein s are thought to ha ve arisen after the co mmon f eatures in the N-terminus. Although co mmon feature s ar e conserve d, dif ferences in amin o aci d composition oc cur, presumably because of di fferent environments of pat hogens, or sy mbiotic relationships with their hosts.

The identified discriminating fe atures were used for cla ssification, and f or elu cidating their interdepe ndence usin g graphical mode ling that comb ined non-linear nonpara metric regression and Bayesian network. Although these techniques are usually used for estimating gene network s from microarray expression data, the combination of them, with featur e selection, wa s a powerf ul metho d for a deeper under standing o f the meaning of t he discriminating features.

This is the fir st study to ex plore discri minating features between pathog enic and symbiotic T3SS ef fector prot eins, usin g a combination of computational stati stical and machine-learning approache s. The mo st-discriminating features, their interde pendence, and the most-discriminating region were determined by these methods. This study will provide a methodological basis for futu re research, and provides important insight about t he functional differences between pathogenic and symbiotic T3SS effectors.

## Referecnes

1.  Cornelis, G. R. (2006) *Nat Rev Microbiol* **4,** 811-25.
2.  Coburn, B., Sekirov, I. & Finlay, B. B. (2007) *Clin Microbiol Rev* **20,** 535-49.
3.  Hernandez, L. D., P ypaert, M., F lavell, R.  A. & Galan,  J. E. (2003) *J Cell Biol* **163,** 1123-31.
4.  Boyle, E. C., Brown, N. F. & Finlay, B. B. (2006) *Cell Microbiol* **8,** 1946-57.
5.  Yoshida, S., Katayama, E., Kuwae, A., Mimuro, H., Suzuki, T. & Sasakawa, C. (2002) *Embo J* **21,** 2923-35.
6.  Beeckman, D. S. & Vanrompay, D. C. (2009) *Curr Issues Mol Biol* **12,** 17-42.
7.  Coombes, B. K. (2009) *Trends Microbiol* **17,** 89-94.
8.  Horn, M., Collingro, A., Schmitz-Esser, S., Beier, C. L., Purk hold, U., Fart mann, B., Brandt, P., Nyakatura, G. J., Droege, M., Frishman, D., Rattei, T., Mewes, H. W. & Wagner, M. (2004) *Science* **304,** 728-30.
9.  Arnold, R., Brandm aier, S., Kleine, F., Tischler, P., Heinz, E., Behrens, S., Niinikoski, A., Mewes, H. W., Horn, M. & Rattei, T. (2009) *PLoS Pathog* **5,** e1000376.
10. Samudrala, R., Heffron, F. & McDermott, J. E. (2009) *PLoS Pathog* **5,** e1000375.
11. Grynberg, M. & Godzik, A. (2009) *PLoS Pathog* **5,** e1000398.
12. Kambara, K., Ardissone, S., Kobayashi, H., Saad, M. M., Schumpp, O., Broughton, W. J. & Deakin, W. J. (2009) *Mol Microbiol* **71,** 92-106.
13. Masson-Boivin, C., Girau d, E.,  Perret, X. &  Batut, J. (2009)  *Trends Micr obiol* **17,** 458-66.
14. Dale, C. & Moran, N. A. (2006) *Cell* **126,** 453-465.
15. Lower, M. & Schneider, G. (2009) *PLoS One* **4,** e5917.
16. Rice, P., Longden, I. & Bleasby, A. (2000) *Trends Genet* **16,** 276-7.
17. Bendtsen, J. D. , Nielsen, H.,  von Heijne,  G. & Brunak,  S. (2004)  *J Mol Biol*  **340,** 783-95.
18. Lepage,    Y. (1971) *Biometrika* **58,** 213-217.
19. Zhu, J. & Hastie,  T. (2001)  *Journal of Computation al and  Graphical S tatistics* **14,** 1081-1088.
20. Cawley, G. C. & Talbot, N. L. (2008) *Machine Learning* **71,** 243 - 264.
21. Konishi, S., Ando, T. & Imoto, S. (2004) *Biometrika* **91,** 27-43.
22. Imoto, S., Goto, T. & Miyano, S. (2002) *Pac Symp Biocomput***,** 175-86.
23. Imoto, S., Suny ong, K., Goto, T., Aburatani, S., Tashiro, K., Kuhara, S. & Miy ano, S. (2002) *Proc IEEE Comput Soc Bioinform Conf* **1,** 219-27.
24. Nabney, I. (2001) *NETLAB: algorithms for pattern recognition* (Springer, London).
25. Murphy, K. P. (2001) *Computing Science and Statistics.* **33,** 331-350.
26. Leray, P. & Francois, O. (2006)    (Laboratoire PSI, Universitè et INSA de Rouen.
27. von Heijne, G. (1985) *J Mol Biol* **184,** 99-105.